

# CSE 444: Database Internals

## Lectures 17 Transactions: Recovery

# Transaction Management

Two parts:

- Concurrency control: ACID
- Recovery from crashes: ACID

We already discussed concurrency control  
You will get to implement locking in lab3

Today, we start recovery  
More details in the book in Chapter 17

# Problem Illustration

```
Client 1:  
START TRANSACTION  
INSERT INTO SmallProduct(name, price)  
SELECT pname, price  
FROM Product  
WHERE price <= 0.99  
  
DELETE Product  
WHERE price <=0.99  
COMMIT
```

Crash !

What do we do now?

# Recovery

From which events below can DBMS recover ?

- Wrong data entry
- Disk failure
- Fire / earthquake / etc.
- **Systems crashes**
  - Software errors
  - Power failures

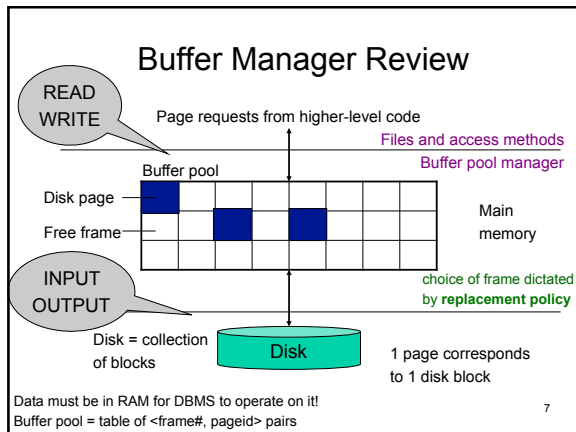
# Recovery

Type of Crash	Prevention
Wrong data entry	Constraints and Data cleaning
Disk crashes	Redundancy: RAID, backup, replica
Fire or other major disaster	Redundancy: Replica far away
<b>System failures</b>	<b>DATABASE RECOVERY</b>

Most frequent

# System Failures

- Each transaction has *internal state*
- When system crashes, internal state is lost
  - Don't know which parts executed and which didn't
  - Need ability to *undo* and *redo*



### Buffer Manager Review

- Enables higher layers of the DBMS to assume that needed data is in main memory
- Caches data in memory. When crash occurs:
  - Problem if committed data was not yet written to disk
  - Problem if uncommitted data was flushed to disk

Magda Balazinska - CSE 444, Spring 2012

### Buffer Manager

- DBMSs build their own buffer manager and don't rely on the OS. Why?
  - Reason 1: Performance
    - DBMS may be able to anticipate access patterns
    - Hence, may also be able to perform prefetching
    - May select better page replacement policy
  - Reason 2: Correctness
    - DBMS needs fine grained control for transactions
    - Needs to force pages to disk for recovery purposes

Magda Balazinska - CSE 444, Spring 2012

### Transactions

- Assumption: db composed of **elements**
  - Usually 1 element = 1 block
  - Can be smaller (=1 record) or larger (=1 relation)
- Assumption: each transaction reads/writes some elements

Magda Balazinska - CSE 444, Spring 2012

### Primitive Operations of Transactions

- READ(X,t)
  - copy element X to transaction local variable t
- WRITE(X,t)
  - copy transaction local variable t to element X
- INPUT(X)
  - read element X to memory buffer
- OUTPUT(X)
  - write element X to disk

Magda Balazinska - CSE 444, Spring 2012

### Example

```

START TRANSACTION
READ(A,t);
t := t*2;
WRITE(A,t);
READ(B,t);
t := t*2;
WRITE(B,t);
COMMIT;

```

Atomicity:  
 BOTH A and B  
 are multiplied by 2

Magda Balazinska - CSE 444, Spring 2012

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)				8	8	
READ(A,t)						
t:=t*2						
WRITE(A,t)						
INPUT(B)						
READ(B,t)						
t:=t*2						
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)						
t:=t*2						
WRITE(A,t)						
INPUT(B)						
READ(B,t)						
t:=t*2						
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)						
INPUT(B)						
READ(B,t)						
t:=t*2						
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)						
READ(B,t)						
t:=t*2						
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)	16	16	8	8	8	
READ(B,t)						
t:=t*2						
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

READ(A,t); t := t\*2; WRITE(A,t);  
READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)						
OUTPUT(A)						
OUTPUT(B)						

Transaction: READ(A,t); t := t\*2; WRITE(A,t);  
 Buffer pool: READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	
OUTPUT(A)						
OUTPUT(B)						


Transaction: READ(A,t); t := t\*2; WRITE(A,t);  
 Buffer pool: READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)						

Transaction: READ(A,t); t := t\*2; WRITE(A,t);  
 Buffer pool: READ(B,t); t := t\*2; WRITE(B,t);

Action	Transaction		Buffer pool		Disk	
	t	Mem A	Mem B	Disk A	Disk B	
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	

Action	t	Mem A	Mem B	Disk A	Disk B
INPUT(A)		8		8	8
READ(A,t)	8	8		8	8
t:=t*2	16	8		8	8
WRITE(A,t)	16	16		8	8
INPUT(B)	16	16	8	8	8
READ(B,t)	8	16	8	8	8
t:=t*2	16	16	8	8	8
WRITE(B,t)	16	16	16	8	8
OUTPUT(A)	16	16	16	16	8
OUTPUT(B)	16	16	16	16	16

 Crash!

Crash occurs after OUTPUT(A), before OUTPUT(B)  
 We lose atomicity

22

### Buffer Manager Policies

- **STEAL or NO-STEAL**
  - Can an update made by an uncommitted transaction overwrite the most recent committed value of a data item on disk?
- **FORCE or NO-FORCE**
  - Should all updates of a transaction be forced to disk before the transaction commits?
- Easiest for recovery: NO-STEAL/FORCE
- Highest performance: STEAL/NO-FORCE

Magda Balazinska - CSE 444, Spring 2012 23

### Solution: Use a Log

- Log = append-only file containing log records
- Note: multiple transactions run concurrently, log records are **interleaved**
- After a system crash, use log to:
  - Redo some transactions that did commit
  - Undo other transactions that did not commit
- Three kinds of logs: undo, redo, undo/redo

Magda Balazinska - CSE 444, Spring 2012 24

## Undo Logging

### Log records

- **<START T>**
  - Transaction T has begun
- **<COMMIT T>**
  - T has committed
- **<ABORT T>**
  - T has aborted
- **<T,X,v>** -- Update record
  - T has updated element X, and its old value was v

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,8>
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
COMMIT						<COMMIT T>

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,8>
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
COMMIT						<COMMIT T>

WHAT DO WE DO ?

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,8>
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
COMMIT						<COMMIT T>

WHAT DO WE DO ?

## After Crash

- In the first example:
  - We UNDO both changes: A=8, B=8
  - The transaction is atomic, since none of its actions have been executed
- In the second example
  - We don't undo anything
  - The transaction is atomic, since both it's actions have been executed

## Undo-Logging Rules

U1: If T modifies X, then <T,X,v> must be written to disk before OUTPUT(X)

U2: If T commits, then OUTPUT(X) must be written to disk before <COMMIT T>

- Hence: OUTPUTs are done early, before the transaction commits

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
INPUT(A)		8		8	8	
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,8>
INPUT(B)	16	16	8	8	8	
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
COMMIT						<COMMIT T>

31

## Recovery with Undo Log

After system's crash, run recovery manager

- Idea 1. Decide for each transaction T whether it is completed or not
  - <START T>.....<COMMIT T>..... = yes
  - <START T>.....<ABORT T>..... = yes
  - <START T>..... = no
- Idea 2. Undo all modifications by incomplete transactions

Magda Balazinska - CSE 444, Spring 2012

32

## Recovery with Undo Log

Recovery manager:

- Read log from the end; cases:
  - <COMMIT T>: mark T as completed
  - <ABORT T>: mark T as completed
  - <T,X,v>: if T is not completed then write X=v to disk else ignore
  - <START T>: ignore

Magda Balazinska - CSE 444, Spring 2012

33

## Recovery with Undo Log

```

...
...
<T6,X6,v6>
...
...
<START T5>
<START T4>
<T1,X1,v1>
<T5,X5,v5>
<T4,X4,v4>
<COMMIT T5>
<T3,X3,v3>
<T2,X2,v2>

```

crash

Question 1 in class:  
Which updates are undone ?

Question 2 in class:  
What happens if there is a second crash, during recovery ?

Question 3 in class:  
How far back do we need to read in the log ?

34

## Recovery with Undo Log

- Note: all undo commands are *idempotent*
  - If we perform them a second time, no harm done
  - E.g. if there is a system crash during recovery, simply restart recovery from scratch

Magda Balazinska - CSE 444, Spring 2012

35

## Recovery with Undo Log

When do we stop reading the log ?

- We cannot stop until we reach the beginning of the log file
- This is impractical

Instead: use checkpointing

Magda Balazinska - CSE 444, Spring 2012

36



## Redo Logging

Log records

- <START T> = transaction T has begun
- <COMMIT T> = T has committed
- <ABORT T> = T has aborted
- <T,X,v> = T has updated element X, and its *new* value is v

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,16>
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,16>
						<COMMIT T>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	

## Redo-Logging Rules

R1: If T modifies X, then both <T,X,v> and <COMMIT T> must be written to disk before OUTPUT(X)

- Hence: OUTPUTs are done *late*

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,16>
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,16>
						<COMMIT T>
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	

## Recovery with Redo Log

After system's crash, run recovery manager

- Step 1. Decide for each transaction T whether it is completed or not
  - <START T>.....<COMMIT T>..... = yes
  - <START T>.....<ABORT T>..... = yes
  - <START T>..... = no
- Step 2. Read log from the beginning, redo all updates of *committed* transactions

## Recovery with Redo Log

```

<START T1>
<T1,X1,v1>
<START T2>
<T2, X2, v2>
<START T3>
<T1,X3,v3>
<COMMIT T2>
<T3,X4,v4>
<T1,X5,v5>
...
...
    
```



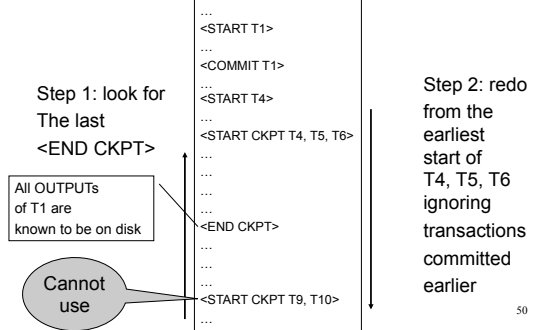
## Nonquiescent Checkpointing

- Write a  $\langle \text{START CKPT}(T_1, \dots, T_k) \rangle$  where  $T_1, \dots, T_k$  are all active transactions
- Flush to disk all blocks of committed transactions (*dirty blocks*), while continuing normal operation
- When all blocks have been written, write  $\langle \text{END CKPT} \rangle$

Magda Balazinska - CSE 444, Spring 2012

49

## Redo Recovery with Nonquiescent Checkpointing



50

## Comparison Undo/Redo

- **Undo logging:**
  - OUTPUT must be done early **Steal/Force**
  - If  $\langle \text{COMMIT } T \rangle$  is seen,  $T$  definitely has written all its data to disk (hence, don't need to redo) – inefficient
- **Redo logging**
  - OUTPUT must be done late **No-Steal/No-Force**
  - If  $\langle \text{COMMIT } T \rangle$  is not seen,  $T$  definitely has not written any of its data to disk (hence there is not dirty data on disk, no need to undo) – inflexible
- Would like more flexibility on when to OUTPUT: **undo/redo logging** (next) **Steal/No-Force**

Magda Balazinska - CSE 444, Spring 2012

51

## Undo/Redo Logging

Log records, only one change

- $\langle T, X, u, v \rangle = T$  has updated element  $X$ , its *old* value was  $u$ , and its *new* value is  $v$

Magda Balazinska - CSE 444, Spring 2012

52

## Undo/Redo-Logging Rule

UR1: If  $T$  modifies  $X$ , then  $\langle T, X, u, v \rangle$  must be written to disk before  $\text{OUTPUT}(X)$

Note: we are free to OUTPUT early or late relative to  $\langle \text{COMMIT } T \rangle$

Magda Balazinska - CSE 444, Spring 2012

53

Action	T	Mem A	Mem B	Disk A	Disk B	Log
						$\langle \text{START } T \rangle$
REAT(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	$\langle T, A, 8, 16 \rangle$
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	$\langle T, B, 8, 16 \rangle$
OUTPUT(A)	16	16	16	16	8	
						$\langle \text{COMMIT } T \rangle$
OUTPUT(B)	16	16	16	16	16	

Can OUTPUT whenever we want: before/after COMMIT<sup>54</sup>

## Recovery with Undo/Redo Log

After system's crash, run recovery manager

- Redo all committed transaction, top-down
- Undo all uncommitted transactions, bottom-up

## Recovery with Undo/Redo Log

